DOCUMENT RESUME

ED 045 734                                    TM 000 325

AUTHOR        Schluck, Gerald J.
TITLE         Miltivariate Sequential Analysis for Education.
INSTITUTION   Florida State Univ., Tallahassee.
NOTE          11p.

EDRS PRICE    EDRS Price MF-$0.25 HC-$0.65
DESCRIPTORS   *Analysis of Variance, Data Analysis, *Hypothesis
              Testing, *Mathematical Applications, Mathematical
              Models, Probability, *Research Methodology,
              Sampling, *Statistical Analysis, Tests of
              Significance
IDENTIFIERS   *Sequential Analysis

ABSTRACT
              Statistical methods that can be applied in the
sequential analysis of multivariate empirical data are provided.
Twenty-three specific formulas for use under varying conditions are
discussed. A historical sketch of sequential analysis since World War
II and a bibliography are included. (AE)

Gerald J. Schluck
Educational Research and Testing
The Florida State University

## MULTIVARIATE SEQUENTIAL ANALYSIS FOR EDUCATION

Given a basic (simple) hypothesis testing situation with null hypothesis,
alternative hypothesis, probability of Type 1 error and probability of Type 11
error, the sequential analysis (Wald) method involves the construction of the
likelihood ratio at each stage of sampling and comparison of the ratio to boundary
values, A and B. If the upper limit is broken with the first sample, the null
hypothesis is rejected. If the lower limit is broken, the alternative hypothesis
is rejected (or the null hypothesis is accepted). If neither limit is broken, a
second sample is selected and the product of the two likelihood ratios is compared
to the boundary values, A and B. The procedure is generally continued until one
of the two limits is broken. Using natural logarithms, the procedure may be repre-
sented as in diagram 1.

Implicit calculations of the upper and lower boundaries involve solution of
the simultaneous equations found in 2. Some computers can handle this
problem but more expedient methods are usually taken. Common estimates for A and
B are found in section 3. Use of these estimates introduce new error probabilities
into the hypothesis testing situation -- it is somewhat comforting to know that the
sum of the new error probabilities is less than or equal to the sum of the desired
error probabilities.

The sequential analysis method for testing a hypothesis necessarily makes the
sample size, say N, a random variable -- N is defined as the size of the sample
when a boundary is first broken. Calculation of the expected value of N from first
conditions involves the summation of the series found in 4. Again, direct calcu-
lations are difficult. Approximations for this mean value exist if the mean
value of the likelihood ratio is not unity. For the particular case in which

population means are being tested, the approximation for the expected value
of the random variable N, given the particular population mean value $\mu$ is given
in equation 5. Equation 6 gives this approximation for the special cases, $\mu_0$ and $\mu_a$.

The power of the test is needed for equation 5. From first conditions, the
power is found by summing the infinite series in equation 7. Often, the power is
approximated using a two step procedure. Under mild conditions, there exists
a value $h(\mu)$ such that the expected value of the likelihood ratio raised to this
power is equal to the constant, 1.0. Equation 8 gives two forms of this statement.
If a non-zero value of $h(\mu)$ exists, the power of the sequential test may be estimated
by the expression found in 9.

I would now like to summarize the above by means of a common example. Assume
the random vector with p elements is distributed normally with known covariance
matrix $\Sigma$. The null hypothesis states the population mean vector is $\mu_0$. The al-
ternative hypothesis states the population mean vector is $\mu_a$. The value z for this
special case is given in section 10 -- a linear combination of the p variates.
For the special case in which p=1, section 11 gives the common expression and
the expression more often used in quality control work.

Continuing with the p-dimensional situation, the approximation for expected
number of vectors which would have to be sampled before a decision is made if the
null hypothesis is assumed to be true is given in equation 12. A similar form
exists if the alternative hypothesis is assumed to be true. For other values of
the population mean vector, this approximation involves the power of the
sequential test. Section 13 gives the details.

Perhaps now is the time for the question, "Why, sequential analysis?" The
procedure was developed by Wald during World War II -- its optimality is expressed
as follows: the expected sample size under either $H_0$ or $H_a$ is usually a fraction
of the fixed sample size test with comparable error probabilities. In the univariate
case for means, the fraction is often $\frac{1}{2}$ or less. Empirical work at The Florida

State University for a limited number of cases with dimensionality two or three

show empirically that the fraction in these multivariate cases if often close to

1/2. Comparable economy in other multivariate cases appears possible.

All of the above comments have implied that one vector is sampled at each

stage of testing procedure. It might be desirable to sample groups of vectors at

each stage of the testing procedure. For the above special case, the only correc-

tion would be a group size factor in the denominator of each term in the variance-

covariance matrix. The expected sample size could be calculated as in equation 14.

The author has taken another approach to sequential analysis in the multivariate

case. Using random normal vectors, the Chi-Square or F (possible Beta) distribu-

tions are used in the likelihood ratio. Sampling is performed with groups of

vectors. The advantages of this procedure are:

 A. The non-centrality parameter is constant for the F test
  corresponding to the Chi-Square test. Relationships between
  the two tests are possible.

 B. All approximations given in the first part of this paper apply
  to this testing procedure.

Certain expected values of random variables are given in section 15, 16, and 17.

The simple hypothesis situation which was introduced before is summarized

in section 18. If the first sample of n vectors yields a value $z_1$ which satisfies

the statement of inequality, a second sample of n vectors is taken and the sum

$z_1 + z_2$ is compared to the boundaries. The complexity of the calculations assume

the experimenter has a programmer available. The expected number of vectors which

would have to be sampled if the null hypothesis is assumed to be true is

approximated by the expression in section 19.

The comparable test with unknown dispersion matrix is summarized in section

20. One comment -- the sample variance-covariance matrix is calculated for each

group of measurements. The expected number of vectors which would have to be

sampled if the null hypothesis is assumed to be true is given, approximately,

by equation 21.

Section 22 gives the basic approach for multivariate analysis of variance.

Section 23 gives the sequential "T" test in the univariate case.

## Closing Comments:

1. It is not necessary to insure that the sample size in each group of vectors is identical if the common approximations are not to be used.

2. No literature exists for optimal group sizes to be used in the sequential analysis.

3. The procedure can be extended to any standard hypothesis testing situation.

Gerald J. Schluck
Educational Research and Testing
The Florida State University

Formulae for
## MULTIVARIATE SEQUENTIAL ANALYSIS FOR EDUCATION

\*\*\*\*\*\*

1.  $f_0(\underset{\sim}{x})$  is the density function of the random vector $\underset{\sim}{x}$ under the null hypothesis, $H_0$.

   $f_a(\underset{\sim}{x})$  is the density function of the random vector $\underset{\sim}{x}$ under the alternative hypothesis, $H_a$.

   $$z_i = \ln \frac{f_a(\underset{\sim}{x_i})}{f_0(\underset{\sim}{x_i})}$$

   (Reject $H_a$)    $\ln B < \sum_{i=1}^{i=k} z_i < \ln A$    (Reject $H_0$)

\*\*\*\*\*\*

2.  $q = P$ (Rejecting $H_0 \mid H_0$);    $\beta = P$ (Rejecting $H_a \mid H_a$)

   $q = P(z_1 \geq \ln A \mid H_0) + P(\ln B < z_1 < \ln A, z_1 + z_2 \geq \ln A \mid H_0) + \ldots$

   $\beta = P(z_1 \leq \ln B \mid H_a) + P(\ln B < z_1 < \ln A, z_1 + z_2 \leq \ln B \mid H_a) + \ldots$

\*\*\*\*\*\*

3.
   $$A \approx \frac{1-\beta}{q} \quad ; \quad B \approx \frac{\beta}{1-q} \ .$$

\*\*\*\*\*\*

4.  For a specified population mean vector $\underset{\sim}{\mu}$ ,

   $$E_{\underset{\sim}{\mu}}(N) = 1 \cdot \left[ 1 - P(\ln B < z_1 < \ln A \mid \underset{\sim}{\mu}) \right] + 2 \cdot \left[ 1 - P(\ln B < z_1 < \ln A, \ln B < z_1 + z_2 < \ln A \mid \underset{\sim}{\mu}) \right] + \ldots$$

\*\*\*\*\*\*

5.  Given $\underset{\sim}{\mu}$, $P$ (Rejecting $H_0 \mid \underset{\sim}{\mu}) = 1 - \beta (\underset{\sim}{\mu})$.

   $$E_{\underset{\sim}{\mu}}(N) \approx \frac{\left[ 1 - \beta (\underset{\sim}{\mu}) \right] \cdot \ln A + \beta (\underset{\sim}{\mu}) \cdot \ln B}{E_{\underset{\sim}{\mu}}(z_i)}$$

\*\*\*\*\*\*

6.  For $\underset{\sim}{\mu} = \underset{\sim}{\mu}_0$ ,

   $$E_{\underset{\sim}{\mu}_0}(N) \approx \frac{q \cdot \ln A + (1 - q) \ln B}{E_{\underset{\sim}{\mu}_0}(z_i)} \ .$$

   For $\underset{\sim}{\mu} = \underset{\sim}{\mu}_a$ ,   $E_{\underset{\sim}{\mu}_a}(N) \approx \frac{(1 - \beta) \ln A + \beta \cdot \ln B}{E_{\underset{\sim}{\mu}_a}(z_i)} \ .$

****** 

7. $1 - \beta(\underset{\sim}{\mu}) = P(z_1 \geq \ln A | \underset{\sim}{\mu}) + P(\ln B < z_1 < \ln A, z_1 + z_2 \doteq \ln A | \underset{\sim}{\mu}) + \ldots$

****** 

8. For population mean $\underset{\sim}{\mu}$, $h(\underset{\sim}{\mu})$ must satisfy

$$E_{\underset{\sim}{\mu}}(e^{h(\underset{\sim}{\mu}) \cdot z_i}) = E_{\underset{\sim}{\mu}}\left[\left(\frac{f_a(\underset{\sim}{x}_i)}{f_o(\underset{\sim}{x}_i)}\right)^{h(\underset{\sim}{\mu})}\right] = 1.$$

****** 

9. For a specified population mean vector $\underset{\sim}{\mu}$ and non-zero $h(\underset{\sim}{\mu})$,

$$1 - \beta(\underset{\sim}{\mu}) \approx \frac{1 - B^{h(\underset{\sim}{\mu})}}{A^{h(\underset{\sim}{\mu})} - B^{h(\underset{\sim}{\mu})}}.$$

****** 

10. $\underset{\sim}{x} \overset{d}{=} N(\underset{\sim}{\mu}, \underset{\sim}{\Sigma})$.

$z_i = \ln f_a(\underset{\sim}{x}_i) - \ln f_o(\underset{\sim}{x}_i) = \underset{\sim}{C}\underset{\sim}{x}_i + \underset{\sim}{D}$,

with $\underset{\sim}{C} = (\underset{\sim}{\mu}_a - \underset{\sim}{\mu}_o)^t \underset{\sim}{\Sigma}^{-1}$ and $\underset{\sim}{D} = \dfrac{\underset{\sim}{\mu}_o^t \underset{\sim}{\Sigma}^{-1} \underset{\sim}{\mu}_o - \underset{\sim}{\mu}_a^t \underset{\sim}{\Sigma}^{-1} \underset{\sim}{\mu}_a}{2}$.

(Reject $H_a$) $\qquad \ln B < \sum_{i=1}^{i=k} (\underset{\sim}{C}\underset{\sim}{x}_i + \underset{\sim}{D}) < \ln A \qquad$ (Reject $H_o$)

****** 

11. $$\ln B < \sum_{i=1}^{i=k}\left[\frac{(\mu_a - \mu_o) \cdot x_i}{\sigma^2} + \frac{\mu_o^2 - \mu_a^2}{2\sigma^2}\right] < \ln A.$$

For $\mu_a > \mu_o$,

$$\left(\ln B\right)\left(\frac{\sigma^2}{\mu_a - \mu_o}\right) + k \cdot \left(\frac{\mu_a + \mu_o}{2}\right) < \sum_{i=1}^{i=k} x_i < \left(\ln A\right) \cdot \left(\frac{\sigma^2}{\mu_a - \mu_o}\right) + k \cdot \left(\frac{\mu_a + \mu_o}{2}\right).$$

****** 

12. $$E_{\underset{\sim}{\mu}_o}(N) \approx \frac{d \cdot \ln A + (1 - d) \ln B}{E_{\underset{\sim}{\mu}_o}(z_i)} = \frac{d \cdot \ln A + (1 - d) \ln B}{\underset{\sim}{C}\underset{\sim}{\mu}_o + \underset{\sim}{D}}$$

with C and D defined in 10 above.

****** 

13. For $\underset{\sim}{x} \overset{d}{=} N(\underset{\sim}{\mu}, \underset{\sim}{\Sigma})$ and $z_i = \underset{\sim}{C}\underset{\sim}{x}_i + \underset{\sim}{D}$,

$$E_{\underset{\sim}{\mu}}(e^{h(\underset{\sim}{\mu})z_i}) = \exp\left(\underset{\sim}{C}\underset{\sim}{\mu} + \underset{\sim}{D} + \frac{h^2(\underset{\sim}{\mu}) \underset{\sim}{C}\underset{\sim}{\Sigma}\underset{\sim}{C}^t}{2}\right) = 1.$$

Certain values of the population mean vector yield a non-zero value $h(\underset{\sim}{\mu})$.

13, cont.

$$h(\underset{\sim}{\mu}) = \frac{-2(\underset{\sim}{C}\,\underset{\sim}{\mu} + D)}{\underset{\sim}{C}\,\underset{\sim}{\Sigma}\,\underset{\sim}{C}^t} \ .$$

For $h(\underset{\sim}{\mu}) \neq 0$,

$$1 - \beta(\underset{\sim}{\mu}) \approx \frac{1 - B^{h(\underset{\sim}{\mu})}}{A^{h(\underset{\sim}{\mu})} - B^{h(\underset{\sim}{\mu})}}$$

$$E_{\underset{\sim}{\mu}}(N) \approx \frac{\left[1 - \beta(\underset{\sim}{\mu})\right] \cdot \ln A + \beta(\underset{\sim}{\mu}) \cdot \ln B}{\underset{\sim}{C}\,\underset{\sim}{\mu} + \underset{\sim}{D}}$$

with $\underset{\sim}{C}$, $\underset{\sim}{D}$, defined in section 10.

******

14.    n vectors are sampled at each stage.

$f_a(\underset{\sim}{\bar{x}}_i)$, $f_o(\underset{\sim}{\bar{x}}_i)$ are density functions of the sample mean vector.

$$z_i = \ln f_a(\underset{\sim}{\bar{x}}_i) - \ln f_o(\underset{\sim}{\bar{x}}_i)$$

$$E_{\underset{\sim}{\mu}}[N] \approx n \cdot \left[\frac{\left[1 - \beta(\underset{\sim}{\mu})\right] \cdot \ln A + \beta(\underset{\sim}{\mu}) \cdot \ln B}{E_{\underset{\sim}{\mu}}(z_i)}\right].$$

******

15.    x follows the non-central chi-square distribution with p degrees of freedom and non-centrality parameter $\delta^2$.

$$E[x] = p + \delta^2$$

$$\text{Var}(x) = 2p + 4\delta^2$$

$$E[\ln x] = \ln 2 + e^{-\frac{\delta^2}{2}} \sum_{i=0}^{\infty} \frac{\left(\frac{\delta^2}{2}\right)^i}{i!} \psi\left(\frac{p}{2} + i\right)$$

with $\psi(a) = \dfrac{\Gamma'(a)}{\Gamma(a)}$, the Digamma (or Psi) function.

******

16.    y follows the non-central F distribution with m and n degrees of freedom and non-centrality parameter $\delta^2$.

$$E[y] = \frac{n}{n-2}\left(1 - \frac{\delta^2}{m}\right) \qquad\qquad n > 2.$$

$$E[\ln y] = \ln\left(\frac{n}{m}\right) - \psi\left(\frac{n}{2}\right) + e^{-\frac{\delta^2}{2}} \sum_{i=0}^{\infty} \frac{\left(\frac{\delta^2}{2}\right)^i}{i!} \psi\left(\frac{m}{2} + i\right).$$

17.    x follows the non-central Beta distribution with a and b degrees of freedom and non-centrality parameter $\delta^2$.

$$E[x] = e^{-\frac{\delta^2}{2}} \cdot a \cdot \sum_{i=0}^{\infty} \frac{(\frac{\delta^2}{2})^i}{i!} \cdot \frac{1}{a+b+i}$$

$$E[\ln x] = \psi(a) - e^{-\frac{\delta^2}{2}} \sum_{i=0}^{\infty} \frac{(\frac{\delta^2}{2})^i}{i!} \cdot \psi(a+b+i) \,.$$

******

18.    n vectors are sampled; $p \geq 2$.

$$z_i = \ln \frac{f\left(\chi^2_{a:i}\right)}{f\left(\chi^2_{o:i}\right)} = \left(\frac{p}{2} - 1\right) \cdot \left[\ln \frac{\chi^2_{a:i}}{\chi^2_{o:i}}\right] - \tfrac{1}{2}\left[\chi^2_{a:i} - \chi^2_{o:i}\right]$$

with    $\chi^2_{a:i} = n(\bar{x}_i - \mu_a)^t \sum^{-1} (\bar{x}_i - \mu_a)$

$\chi^2_{o:i} = n(\bar{x}_i - \mu_o)^t \sum^{-1} (\bar{x}_i - \mu_o)$.

Inversion of a matrix can be avoided if the following relationship is used:

$$a^t \sum^{-1} a = \frac{|\sum + aa^t|}{|\sum|} - 1 \,.$$

******

19.

$$E_{\mu_o}[N] \approx n \cdot \left[\frac{\alpha \cdot \ln A + (1 - \alpha) \ln B}{E_{\mu_o}[z_i]}\right]$$

$$E_{\mu_o}(z_i) = \left(\frac{p}{2} - 1\right) \left[e^{-\frac{\delta^2}{2}} \cdot \sum \frac{(\frac{\delta^2}{2})^i}{i!} \psi\left(\frac{p}{2} + i\right) - \psi\left(\frac{p}{2}\right)\right] - \tfrac{1}{2}\delta^2$$

with $\delta^2 = n(\mu_a - \mu_o)^t \sum^{-1} (\mu_a - \mu_o)$.

******

20.      n vectors are sampled;    $p \geq 2$.

$$z_i = \ln \frac{g\left(F_{a:i}\right)}{g\left(F_{o:i}\right)} = \left(\frac{p}{2} - 1\right)\left[\ln \frac{F_{a:i}}{F_{o:i}}\right] + \frac{n}{2}\left[\ln \frac{\left(1 + \frac{p}{n-p} F_{o:i}\right)}{\left(1 + \frac{p}{n-p} F_{a:i}\right)}\right]$$

with   $F_{a:i} = \frac{(n-p)}{p(n-1)} \cdot n(\tilde{x}_i - \mu_a) S_i^{-1} (\tilde{x}_i - \mu_a)$

$$F_{o:i} = \frac{(n-p)}{p(n-1)} \cdot n(\tilde{x}_i - \mu_o) S_i^{-1} (\tilde{x}_i - \mu_o).$$

\*\*\*\*\*\*

21.

$$E_{\mu_o}[N] \approx n \cdot \left[\frac{\alpha \cdot \ln A + (1 - \alpha)\ln B}{E_{\mu_o}[z_i]}\right]$$

$$E_{\mu_o}(z_i) = \left(\frac{p}{2} - 1\right) \cdot \left[e^{-\frac{\delta^2}{2}} \sum_{\lambda = 0}^{\infty} \frac{\left(\frac{\delta^2}{2}\right)^\lambda}{i!} \psi\left(\frac{p}{2} + i\right) - \psi\left(\frac{p}{2}\right)\right]$$

$$- \frac{n}{2} \cdot \left[e^{-\frac{\delta^2}{2}} \sum_{\lambda = 0} \frac{\left(\frac{\delta^2}{2}\right)^\lambda}{i!} \psi\left(\frac{n}{2} + i\right) - \psi\left(\frac{n}{2}\right)\right]$$

\*\*\*\*\*\*

22.      Random samples of n vectors are taken from 3 multivariate normal populations.

$$x_1 \overset{d}{=} N(\mu^{(1)}, \Sigma), \qquad x_2 \overset{d}{=} N(\mu^{(2)}, \Sigma) \quad \text{and} \quad x_3 \overset{d}{=} N(\mu^{(3)}, \Sigma).$$

The hypothesis of equal population mean vectors becomes

$$H_o: \qquad C \cdot \begin{bmatrix} \mu^{(1)} \\ \mu^{(2)} \\ \mu^{(3)} \end{bmatrix} = 0 \; ; \qquad \begin{array}{c} C \text{ is a contrast matrix.} \\ (2p) \times (3p) \end{array}$$

$$C = \begin{bmatrix} I & -I & 0 \\ I & I & -2I \end{bmatrix} \qquad \text{is a common selection.}$$

22. cont.     The alternative hypothesis may be:

$$H_a: \quad \underset{\sim}{C} \begin{bmatrix} \underset{\sim}{\mu}^{(1)} \\ \underset{\sim}{\mu}^{(2)} \\ \underset{\sim}{\mu}^{(3)} \end{bmatrix} = \underset{\sim}{a}_{(2p \times 1)} = \begin{bmatrix} \underset{\sim}{a}_1 \\ \underset{\sim}{a}_2 \end{bmatrix} .$$

Sequential procedures follow immediately with non-centrality parameter in both the known dispension matrix problem and unknown dispension matrix problem equal to

$$\delta 2 = n \left[ \frac{\underset{\sim}{a}_1^t \underset{\sim}{\Sigma}^{-1} \underset{\sim}{a}_1}{2} + \frac{\underset{\sim}{a}_2^t \underset{\sim}{\Sigma}^{-1} \underset{\sim}{a}_2}{6} \right] .$$

\*\*\*\*\*\*

23.   Univariate "t".

n measurements are samples from the normal population.

$$H_o: \mu = \mu_o \; ; \qquad H_a: \mu = \mu_a .$$

$$F_{a:i} = n(\bar{x}_i - \mu_a)^t S^{-1} (\bar{x}_i - \mu_a) = n \frac{(\bar{x}_i - \mu_a)^2}{s_i^2}$$

$$F_{o:i} = n \frac{(\bar{x}_i - \mu_o)^2}{s_i^2} .$$

$$z_i = \frac{n}{2} \ln \left[ \frac{1 + \frac{1}{n-1} F_{o:i}}{1 + \frac{1}{n-1} F_{a:i}} \right]$$

It can be shown that the average sampling number in this situation is greater than the average sampling number in the corresponding known-variance test.

\*\*\*\*\*\*
\*\*\*\*\*\*

## Historical Sketch of Sequential Analysis

The initial breakthrough in sequential analysis was made by Wald (1945, 1947). Cox (1952) and Girshick (1946) extended the basic technique to variance tests and exponential family tests. Johnson (1953, 1954) worked extensively with analysis of variance. Papers on multivariate sequential analysis were written by Jackson and Bradley (1961A, 1961B). Most of the more recent literature deals with arbitrary stopping rules---see Myers, Schneiderman and Armitage (1966). Theoretical work with non-central statistics presented in this paper depend heavily upon the Mellin Transform as presented by Epstein (1948).

## Bibliography

Cox, D. R. Sequential tests for composite hypotheses. Proceedings of Cambridge Philosophical Society, 1952, 48, 290-299.

Epstein, B. Some applications of the Mellin Transform in statistics. Annals of Mathematical Statistics, 1948, 19, 370-379.

Girshick, M. A. Contributions to the theory of sequential analysis, II, III. Annals of Mathematical Statistics, 1946, 17, 282-298.

Jackson, J. E. and Bradley, R. A. Sequential $\chi^2$ and $T^2$ tests. Annals of Mathematical Statistics, 1961, 32, 1063-1077.

Jackson, J. E. and Bradley, R. A. Sequential $\chi^2$ and $T^2$ tests and their application to an acceptance sampling problem. Technometrics, 1961, 3, 519-534.

Johnson, N. L. Some notes on the application of sequential methods in the analysis of variance. Annals of Mathematical Statistics, 1953, 24, 614-623.

Johnson, N. L. Sequential procedures in certain component of variance problems. Annals of Mathematical Statistics, 1954, 25, 357-366.

Myers, M. H. and Schneiderman, M. A. and Armitage, P. Boundaries for closed (wedge) sequential t-test plans. Biometrika, 1966, 53, 431-437.

Wald, A. Sequential tests of statistical hypotheses. Annals of Mathematical Statistics, 1945, 16, 117-186.

Wald, A. Sequential Analysis. New York: Wiley, 1947.